

1 **A 2-Gene Host Signature for Improved Accuracy of COVID-19 Diagnosis Agnostic to Viral**
2 **Variants**

3
4 Jack Albright^{1,*}, Eran Mick^{1,2,3,*}, Estella Sanchez-Guerrero^{2,*}, Jack Kamm^{1,+}, Anthea Mitchell^{1,4},
5 Angela M. Detweiler¹, Norma Neff¹, Alexandra Tsitsiklis², Paula Hayakawa Serpa²,
6 Kalani Ratnasiri¹, Diane Havlir⁵, Amy Kistler¹, Joseph L. DeRisi^{1,4}, Angela Oliveira Pisco¹,
7 Charles R. Langelier^{1,2,‡}

8 * Equal contribution

9
10 ¹ Chan Zuckerberg Biohub, San Francisco, CA, USA.

11 ² Division of Infectious Diseases, Department of Medicine, University of California San Francisco,
12 San Francisco, CA, USA.

13 ³ Division of Pulmonary and Critical Care Medicine, Department of Medicine, University of
14 California San Francisco, San Francisco, CA, USA.

15 ⁴ Department of Biochemistry and Biophysics, University of California San Francisco, San
16 Francisco, CA, USA.

17 ⁵ Division of HIV, Infectious Diseases and Global Medicine, Department of Medicine, University
18 of California San Francisco, San Francisco, CA, USA.

19
20 + Present address: Genentech, Inc., South San Francisco, CA, USA

21 ‡ Correspondence: chaz.langelier@ucsf.edu

22 **Abstract**

23 The continued emergence of SARS-CoV-2 variants is one of several factors that may
24 cause false negative viral PCR test results. Such tests are also susceptible to false positive results
25 due to trace contamination from high viral titer samples. Host immune response markers provide
26 an orthogonal indication of infection that can mitigate these concerns when combined with direct
27 viral detection. Here, we leverage nasopharyngeal swab RNA-seq data from patients with COVID-
28 19, other viral acute respiratory illnesses and non-viral conditions (n=318) to develop support
29 vector machine classifiers that rely on a parsimonious 2-gene host signature to predict COVID-
30 19. Optimal classifiers achieve an area under the receiver operating characteristic curve (AUC)
31 greater than 0.9 when evaluated on an independent RNA-seq cohort (n=553). We show that a
32 classifier relying on a single interferon-stimulated gene, such as *IFI6* or *IFI44*, measured in RT-
33 qPCR assays (n=144) achieves AUC values as high as 0.88. Addition of a second gene, such as
34 *GBP5*, significantly improves the specificity compared to other respiratory viruses. The
35 performance of a clinically practical 2-gene RT-qPCR classifier is robust across common SARS-
36 CoV-2 variants, including Omicron, and is unaffected by cross-contamination, demonstrating its
37 utility for improving accuracy of COVID-19 diagnostics.

38 Introduction

39 The COVID-19 pandemic has inflicted unprecedented human health consequences, with
40 millions of deaths reported worldwide since December 2019¹. Testing is a cornerstone of pandemic
41 management, yet existing assays suffer from accuracy limitations. Even the gold-standard testing
42 modality of nasopharyngeal (NP) swab RT-PCR returns falsely negative in a substantial proportion
43 of cases²⁻⁴ and may fail to detect SARS-CoV-2 variants with mutations at primer target sites⁵⁻⁷.
44 False positive tests due to sample cross-contamination in the laboratory are also a significant
45 complication^{8,9} as they can lead to costly contact tracing efforts and the unnecessary isolation of
46 uninfected individuals, including essential workers.

47 Measuring the host immune response offers a complementary approach to direct
48 detection of the SARS-CoV-2 virus and holds potential for overcoming the limitations of existing
49 COVID-19 diagnostics. RNA-sequencing (RNA-seq) studies of NP swabs and blood have
50 demonstrated that COVID-19 elicits a unique host transcriptional response compared with non-
51 viral and other viral acute respiratory illnesses (ARIs)¹⁰⁻¹². A host gene expression signature of
52 COVID-19, when utilized in combination with molecular detection of SARS-CoV-2, can serve as
53 a fallback to identify suspected false negative or false positive results of traditional viral PCR tests,
54 thus improving overall diagnostic reliability.

55 Recent studies have employed machine learning on RNA-seq data from NP swabs to
56 develop proof-of-concept, host-based COVID-19 diagnostic classifiers that rely on a relatively
57 large number of genes^{10,13}. While highly promising, these classifiers have yet to undergo validation
58 in external cohorts. Furthermore, RNA-seq is not widely available in clinical settings and thus the
59 immediate practical utility of RNA-seq classifiers is limited.

60 Here, we address these gaps by developing 2-gene host signatures that could practically
61 be incorporated into an RT-qPCR (qPCR) assay alongside a control gene and a viral target. We
62 leverage NP swab RNA-seq data from two large patient cohorts to derive and validate top
63 performing support vector machine (SVM) binary classifiers that use 2 host genes to predict

64 COVID-19 status. We then refine these signatures for use in qPCR assays and confirm their
65 prediction performance from qPCR data on an independent sample cohort.

66 The optimal 2-gene signatures combine an interferon-stimulated gene (ISG) that is
67 strongly induced in COVID-19, such as *IFI6*, *IFI44* or *IFI44L*, with another immune response gene
68 that is more strongly induced in other viral ARIs, such as *GBP5* or *CCL3*. Finally, we demonstrate
69 that such a host classifier is robust across SARS-CoV-2 variants, including those that can yield a
70 false negative viral PCR result, and is unaffected by laboratory cross-contamination that can yield
71 a false positive viral PCR result.

72

73 **Results**

74

75 Performance of 2-gene combinations for a COVID-19 host classifier from NP swab RNA-seq

76 We previously developed multi-gene host classifiers for COVID-19 using RNA-seq data
77 from NP swabs of patients tested for COVID-19 at the University of California, San Francisco
78 (UCSF) who were diagnosed with either COVID-19, other viral ARIs or non-viral ARIs¹⁰. In the
79 present work, we sought to develop a parsimonious 2-gene signature that could practically be
80 incorporated into a PCR test alongside a control gene and a viral target. We began by identifying
81 top performing 2-gene candidates in our RNA-seq cohort after supplementing it with additional
82 samples collected in the intervening time. The full UCSF cohort included n=318 patients, of whom
83 90 had COVID-19 (with viral load equivalent to PCR $C_t < 30$), 59 had other viral ARIs (mostly
84 rhinovirus and influenza), and 169 had non-viral conditions (**Supp. Table 1; Supp. Data File 1**).

85 The UCSF samples were split into a training set (70%) and a testing set (30%), with
86 stratification to ensure each one contained similar proportions of samples with and without
87 COVID-19. We then applied a greedy selection algorithm to identify 2-gene combinations that
88 best predicted COVID-19 status. The performance metric was the area under the receiver
89 operating characteristic curve (AUC) of a support vector machine (SVM) binary classifier that

90 used the selected genes as features, calculated using 5-fold cross-validation within the training
91 set (**Figure 1a**). Thus, a first gene was selected to maximize the AUC it achieved, and a second
92 gene was selected to maximize the AUC when combined with the first gene. **Table 1a** lists nine
93 combinations composed of each of the three best ‘first’ genes and their respective three best
94 ‘second’ genes. The ‘first’ genes in the top combinations were the interferon-stimulated genes
95 (ISGs) *IFI6*, *IFI44L* and *HERC6*, which we previously showed are strongly induced in COVID-
96 19¹⁰. Most of the ‘second’ genes were also related to immune and inflammatory processes.

97 The performance of the nine 2-gene combinations on previously unseen data was
98 estimated by: i) 10,000 rounds of 5-fold cross-validation within the training set, ii) 10,000 rounds
99 of 5-fold cross-validation within the testing set, or iii) training on the training set and prediction on
100 the testing set (**Table 1a**). Using the third approach, we observed AUC values as high as 0.93
101 (**Figure 1b**). We further validated the classifiers using an external, independently generated and
102 quantified NP swab RNA-seq dataset from a cohort of n=553 patients in New York (166 with
103 COVID-19, 79 with other viral ARIs, 308 with non-viral conditions)¹² (**Supp. Table 1; Supp. Data**
104 **File 1**). The 2-gene combinations achieved comparable performance on the external dataset
105 (**Table 1b**). The best performing combinations were *IFI44L+GBP5* (AUC 0.919) and *IFI6+GBP5*
106 (AUC 0.91), when the classifier was trained on the UCSF 70% training set. These results
107 demonstrate that 2-gene classifier models are feasible, stable, and generalizable.

108

109 Optimization of 2-gene combinations for incorporation into an RT-qPCR assay

110 We noted that the ‘first’ and the ‘second’ genes in the 2-gene combinations performed
111 distinct roles. The former was sufficient to distinguish COVID-19 from non-viral ARIs while the
112 latter helped reinforce the distinction between COVID-19 and other viral ARIs. Considering *IFI6*
113 and *GBP5* as an example, *IFI6* alone almost completely separated the COVID-19 and non-viral
114 samples (**Figure 1c**). However, some of the other viral ARI samples showed equivalent levels of
115 *IFI6* expression. Adding *GBP5* to the model allowed for improved separation, as expression of

116 this ISG was typically higher in other viral ARIs (**Figure 1c**). Given this pattern, and because our
117 ultimate goal was a qPCR assay in which small effect sizes are more difficult to discern, we refined
118 our candidate genes by also considering the expression fold-change between COVID-19 and the
119 two other patient groups.

120 We first plotted the AUC of SVM classifiers relying on each individual gene against the
121 fold-change of that gene between the COVID-19 and non-viral samples, where both measures
122 were averaged between the full UCSF cohort and the New York cohort (**Figure 1d; Supp. Data**
123 **File 2**). As expected, several ISGs exhibited equivalently robust predictive value as well as
124 substantial fold-changes ($\log_2FC \sim 2-4$) that should be readily detectable by qPCR. We then
125 plotted the AUC of classifiers that used the ISG *IFI6* in combination with every possible ‘second’
126 gene, against the fold-change of the ‘second’ gene between COVID-19 and other viral ARIs
127 (**Figure 1e; Supp. Data File 2**). This revealed candidate genes with somewhat smaller fold-
128 changes ($\log_2FC \sim 1.5-2$) that should still be detectable by qPCR. These candidate genes only
129 partly overlapped with the ‘second’ genes selected by the greedy algorithm, which did not
130 explicitly consider fold-change.

131

132 RT-qPCR validation of host genes to differentiate COVID-19 from other ARIs

133 We chose four ‘first’ ISGs based on their predictive value and fold-change. We then
134 measured the expression of these ISGs, relative to the reference gene *RPP30*, using qPCR in
135 swabs from a new cohort of patients with (n=72) or without (n=72) COVID-19. Because these
136 swabs were not sequenced, we could not definitively assign those without COVID-19 as either
137 non-viral or other viral cases. However, the low prevalence of other viral ARIs during the
138 timeframe of sample collection, due to the public health measures implemented for COVID-19¹⁴,
139 suggested they were mostly non-viral. All four genes were able to clearly separate the majority of
140 samples with or without COVID-19 in the qPCR data (**Figure 2a**). SVM classifiers relying on single

141 ISGs achieved mean AUC values as high as 0.88 in predicting COVID-19 status from the qPCR
142 data (**Figure 2b**), on par with their prediction performance from RNA-seq (**Figure 1d**).

143 To explicitly test the ability to separate COVID-19 from other respiratory viruses using
144 qPCR data, we chose four ‘second’ genes and measured their expression in the COVID-19
145 samples described above (n=72) as compared to a subset of our original, sequenced other viral
146 samples (n=17). Consistent with the RNA-seq data, the expression of three of the genes (*GBP5*,
147 *CD274* and *CCL3*) was significantly higher on average in the other viral samples, and expression
148 of the fourth gene (*GSTA2*) was significantly higher in the COVID-19 samples (**Figure 2c**).
149 Classifiers that combined one of these genes with a ‘first’ gene achieved near perfect separation
150 of the COVID-19 and other viral samples (**Figure 2d**). This performance is likely overly optimistic,
151 due in part to the relatively small size of the other virus group in the qPCR data, but it is overall
152 consistent with the performance observed in the larger RNA-seq datasets. These results
153 demonstrate that 2-gene signatures can successfully predict COVID-19 status from qPCR data.

154

155 Host signatures are robust to SARS-CoV-2 variants and laboratory cross-contamination

156 We next assessed whether a 2-gene host classifier was robust across SARS-CoV-2
157 variants, which could conceivably yield an altered host response and/or harbor mutations that
158 disrupt primer target sites and lead to false negative viral PCR tests^{5,7,15}. We performed qPCR for
159 the genes *IFI6* and *GBP5* on samples with the Omicron variant (n=3), which causes S-gene target
160 dropout in certain viral PCR assays; the N-gene variant (n=4), which causes N-gene target
161 dropout¹⁵; and on samples with the Delta variant (n=7). An SVM classifier trained on the qPCR
162 results of the samples with and without COVID-19, described above, predicted COVID-19 with
163 high likelihood in all variant samples (**Figure 2e**), demonstrating the utility of a host signature as
164 a complement to viral PCR.

165 On the other hand, false positive viral PCR tests frequently result from trace cross-
166 contamination of samples with high viral titers into negative specimens processed

167 contemporaneously in the laboratory⁹. To examine whether the *IFI6+GBP5* host classifier would
168 be affected in such cross-contamination events, we spiked extracted NP swab RNA from a
169 sample with very high SARS-CoV-2 viral load ($C_t \approx 12$) into $n=7$ COVID-19 negative swab
170 specimens at a dilution of $1:10^5$, which would be expected to yield a positive viral PCR with $C_t < 30$.
171 The probability of COVID-19 estimated by the host classifier was not significantly affected in the
172 simulated false-positive specimens (**Figure 2f**).

173

174 **Discussion**

175 We leveraged multiple cohorts – encompassing over 1,000 patients with COVID-19, other
176 viral ARIs and non-viral conditions – to develop and validate 2-gene host-based COVID-19
177 diagnostic classifiers that could be practically incorporated into clinical PCR assays in
178 combination with a control gene and a viral target. We found that the host classifier enabled
179 reliable identification of COVID-19 even in the face of SARS-CoV-2 variants that cause false
180 negative viral PCR tests and remained unaffected by simulated laboratory cross-contamination
181 that can cause false positive viral PCR tests.

182 Given the inevitable continued emergence of SARS-CoV-2 variants, which may disrupt
183 primer target sites, assays capable of detecting infection regardless of viral sequence are
184 essential to avoid adverse outcomes owing to infected individuals going unrecognized in
185 congregate settings, such as hospitals or nursing homes. The adverse effects of false-positive
186 tests are also non-trivial. The positive predictive value of highly specific viral PCR assays
187 diminishes for asymptomatic individuals undergoing continual surveillance testing in low
188 prevalence settings⁹. False positive results then become more likely, leading to unnecessary
189 isolation and quarantine, depletion of essential personnel, and unwarranted contact tracing.

190 Our study has some limitations. While our findings provide a framework for the rapid
191 clinical translation of a host-based COVID-19 diagnostic, a randomized controlled trial of our
192 assay will be needed to firmly establish its clinical utility. Our results suggest that addition of host

193 targets is likely to improve diagnostic accuracy, however, a prospective assessment using
194 clinically confirmed false-positive and false-negative viral tests is needed. Moreover, our classifier
195 models were trained and tested on cohorts with particular characteristics, including the balance
196 between COVID-19, other viral and non-viral samples; the mix of other respiratory viruses
197 represented; and within the COVID-19 group, the distributions of viral load and of time since onset
198 of infection. All these variables no doubt affect classifier performance and will vary in reality with
199 time and place. However, the fact that our classifiers translated so well across diverse real-world
200 cohorts argues that they are quite robust to these issues.

201 While we did not explicitly explore it here, our results suggest that parsimonious host
202 classifiers could serve not only as a COVID-19 diagnostic but also as a pan-respiratory virus
203 surveillance tool. Even prior to the COVID-19 pandemic, viral lower respiratory tract infections
204 were a leading cause of disease and death¹⁶, and many respiratory viral infections go undetected,
205 leading to preventable transmission and unnecessary antibiotic treatment¹⁷. Since our classifiers
206 rely heavily on ISGs, and type I interferon signaling is a biologically conserved mechanism, these
207 genes could be used in future work as the basis for a diagnostic that identifies respiratory viruses
208 more generally. Such a diagnostic could have considerable value as a screening tool in hospitals,
209 nursing homes or other congregate settings with potential for adverse consequences from
210 unrecognized respiratory viral transmission.

211 **Materials and Methods**

212 RNA-seq cohorts and data pre-processing

214 The UCSF cohort used to develop the RNA-seq classifiers was initially described in our
215 study applying metagenomic sequencing to NP swabs from adult patients tested for COVID-19
216 by RT-PCR, according to the published methods and under UCSF IRB #17-24056¹⁰. In brief,
217 samples were assigned to one of three viral status groups: 1) samples with a positive clinical RT-
218 PCR test for SARS-CoV-2 were assigned to the “COVID-19” group, 2) samples with another

219 pathogenic respiratory virus detected by the ID-Seq pipeline¹⁸ in the metagenomic sequencing
220 data were assigned to the “other virus” group, and 3) remaining samples were assigned to the
221 “no virus” group. In the present work, we supplemented the samples reported in our original study
222 with additional swabs collected, sequenced and analyzed in the same manner.

223 We wished to retain for classifier development COVID-19 samples with likely active
224 infection (culturable virus), which several studies have related to viral PCR $C_t < 30$ ^{19–21}. Because
225 not all C_t values were available, we relied on the relationship between viral reads-per-million (rpM)
226 in the sequencing data and PCR C_t that we previously reported¹⁰: $\log_2(\text{rpM}) = 31.9753 - 0.9167 * C_t$.
227 Metadata for the UCSF samples is provided in **Supp. Data File 1**.

228 We pseudo-aligned the UCSF samples with kallisto²² (v. 0.46.1), using the bias correction
229 setting, against an index consisting of all transcripts associated with human protein coding genes
230 (ENSEMBL v. 99), cytosolic and mitochondrial ribosomal RNA sequences, and the sequences of
231 ERCC RNA standards. Samples retained in the dataset had at least 400,000 estimated counts
232 associated with transcripts of protein coding genes. Gene-level counts were generated from the
233 kallisto transcript abundance estimates using the R package tximport²³ (v. 1.14) with the
234 scaledTPM method. Genes were retained if they had at least 10 counts in at least 20% of samples.

235 The New York cohort used to validate the RNA-seq SVM classifiers on an external dataset
236 was previously published¹². Samples in this cohort were also categorized into the three viral status
237 groups described above based on a combination of RT-PCR and metagenomic sequencing.
238 Because we did not have access to the underlying sequencing data, we used the gene counts
239 originally generated by the authors using STAR alignment and the R function featureCounts. We
240 excluded samples with less than five million total counts as well as samples that had discordant
241 COVID-19 test results between two assays, but did not filter based on viral load. Genes were
242 retained if they had at least 32 counts in at least 10% of samples. Metadata for the New York
243 samples is provided in **Supp. Data File 1**.

244 For each RNA-seq cohort, gene counts were subjected to the variance stabilizing
245 transformation (VST) from the R package DESeq2 (v. 1.26.0) and the transformed values were
246 then standardized (centered and scaled) to yield the final input features.

247 RNA-seq SVM classifier development and validation

248 SVM learning was implemented in scikit-learn (<https://scikit-learn.org>) using the
249 `sklearn.svm.SVC` class function with default parameters and probabilistic output.

250 The UCSF cohort was split into a training set (70%) and a testing set (30%), with
251 stratification to ensure each set contained a similar proportion of samples with and without
252 COVID-19. For the greedy feature selection, performance of a binary SVM classifier for predicting
253 COVID-19 status relying on each single feature (gene) was evaluated by running 5-fold cross-
254 validation within the training set and calculating the average AUC across the folds. The three best-
255 performing 'first' genes were then selected. To extend these 'first' genes to 2-gene combinations,
256 another round of the algorithm was performed, picking the three best-performing 'second' genes
257 when combined with each of the 'first' genes.

258 In order to rigorously assess the performance of the SVM 2-gene models, we employed
259 three approaches: (1) running 10,000 rounds of 5-fold cross-validation on the UCSF 70% training
260 set and calculating the average AUC and standard deviation, (2) running 10,000 rounds of 5-fold
261 cross-validation on the UCSF 30% testing set and calculating the average AUC and standard
262 deviation, and (3) training each model on the UCSF 70% training set and testing it on the 30%
263 testing set to generate an AUC score (**Table 1a**). We then validated the 2-gene models on the
264 external New York cohort, using two approaches: (1) running 10,000 rounds of 5-fold cross-
265 validation on the New York cohort and calculating the average AUC and standard deviation, and
266 (2) training each model on the UCSF 70% training set and testing it on the New York cohort to
267 generate an AUC score (**Table 1b**).

268 The AUC scores of single gene or 2-gene SVM classifiers displayed in **Figure 1d,e** were
269 calculated by 5-fold cross validation using all the samples in each RNA-seq cohort.

270 RNA-seq differential expression

271 Gene expression fold-changes in each RNA-seq cohort between the COVID-19 and non-
272 viral samples (**Figure 1d**) and between the COVID-19 and other viral samples (**Figure 1e**) were
273 calculated with the R package limma (v. 3.42), using quantile normalization and the voom method.

274 RT-qPCR of host genes

275 RNA was reverse transcribed using the High-Capacity cDNA Reverse Transcription Kit
276 (Applied Biosystems), according to the manufacturer's protocol, and analyzed by qPCR in a Bio-
277 Rad CFX384 thermocycler (BioRad) using Taqman Fast Advanced Master Mix (Applied
278 Biosystems) and Taqman Gene Expression Assays (Applied Biosystems), according to the
279 manufacturer's protocol. Assay IDs for each gene are provided in **Supp. Table 2**. ΔC_t values were
280 calculated with respect to the reference gene *RPP30* (also known as *RNASEP2*), the standard
281 host control gene used in many viral PCR tests. ΔC_t values are provided in **Supp. Data File 3**.

282 qPCR SVM classifier development and validation

283 The input features for qPCR-based SVM COVID-19 diagnostic classifiers were
284 standardized (centered and scaled) ΔC_t values. Standardization was performed separately in
285 each analyzed sample set using the mean and standard deviation of the training samples. In the
286 context of cross-validation, this was done for each fold using the appropriate training samples.
287 Performance of SVM classifiers to distinguish between the samples with (n=72) and without
288 (n=72) COVID-19 (**Figure 2b**), or to distinguish between the samples with COVID-19 and other
289 viral ARIs (n=17) (**Figure 2d**), was assessed by 5-fold cross-validation.

290 The *IFI6+GBP5* classifier, which was used to predict the COVID-19 status of variant
291 samples (**Figure 2e**) and of samples that had been purposely contaminated with 1:10⁵ dilution
292 from a high SARS-CoV-2 viral load sample (**Figure 2f**), was trained on the set of samples with

293 and without COVID-19 described above. Because the variant and contamination samples on
294 which we performed prediction were assayed in separate experiments subsequent to the
295 generation of the training dataset, they were always processed alongside n=6-7 COVID-19
296 negative controls from the original training dataset. The median ΔC_t difference observed for these
297 control samples between the training dataset and the prediction experiment in which they were
298 re-run was applied to all the samples in the respective experiment in order to account for
299 systematic shifts.

300 **References**

- 301 1. World Health Organization. WHO COVID-19 Dashboard. <https://covid19.who.int> (2021).
- 302 2. Arevalo-Rodriguez, I. *et al.* False-negative results of initial RT-PCR assays for COVID-19: A
303 systematic review. *PLoS ONE* **15**, e0242958 (2020).
- 304 3. Long, D. R. *et al.* Occurrence and Timing of Subsequent Severe Acute Respiratory Syndrome
305 Coronavirus 2 Reverse-transcription Polymerase Chain Reaction Positivity Among Initially
306 Negative Patients. *Clinical Infectious Diseases* **72**, 323–326 (2021).
- 307 4. Kanji, J. N. *et al.* False negative rate of COVID-19 PCR testing: a discordant testing analysis.
308 *Virology* **18**, 13 (2021).
- 309 5. van Dorp, L. *et al.* Emergence of genomic diversity and recurrent mutations in SARS-CoV-2.
310 *Infection, Genetics and Evolution* **83**, 104351 (2020).
- 311 6. Galloway, S. E. *et al.* Emergence of SARS-CoV-2 B.1.1.7 Lineage — United States,
312 December 29, 2020–January 12, 2021. *MMWR Morb. Mortal. Wkly. Rep.* **70**, 95–99 (2021).
- 313 7. U.S. Food and Drug Administration. Genetic Variants of SARS-CoV-2 May Lead to False
314 Negative Results with Molecular Tests for Detection of SARS-CoV-2 - Letter to Clinical
315 Laboratory Staff and Health Care Providers. [https://www.fda.gov/medical-devices/letters-](https://www.fda.gov/medical-devices/letters-health-care-providers/genetic-variants-sars-cov-2-may-lead-false-negative-results-molecular-tests-detection-sars-cov-2)
316 [health-care-providers/genetic-variants-sars-cov-2-may-lead-false-negative-results-](https://www.fda.gov/medical-devices/letters-health-care-providers/genetic-variants-sars-cov-2-may-lead-false-negative-results-molecular-tests-detection-sars-cov-2)
317 [molecular-tests-detection-sars-cov-2](https://www.fda.gov/medical-devices/letters-health-care-providers/genetic-variants-sars-cov-2-may-lead-false-negative-results-molecular-tests-detection-sars-cov-2) (2021).
- 318 8. Surkova, E., Nikolayevskyy, V. & Drobniewski, F. False-positive COVID-19 results: hidden
319 problems and costs. *The Lancet Respiratory Medicine* **8**, 1167–1168 (2020).
- 320 9. Healy, B., Khan, A., Metezai, H., Blyth, I. & Asad, H. The impact of false positive COVID-19
321 results in an area of low prevalence. *Clin Med* **21**, e54–e56 (2021).
- 322 10. Mick, E. *et al.* Upper airway gene expression reveals suppressed immune responses to
323 SARS-CoV-2 compared with other respiratory viruses. *Nature Communications* **11**, (2020).
- 324 11. McClain, M. T. *et al.* Dysregulated transcriptional responses to SARS-CoV-2 in the periphery.
325 *Nat Commun* **12**, 1079 (2021).
- 326 12. Butler, D. *et al.* Shotgun transcriptome, spatial omics, and isothermal profiling of SARS-CoV-
327 2 infection reveals unique host responses, viral diversification, and drug interactions. *Nat*
328 *Commun* **12**, 1660 (2021).

- 329 13. Ng, D. L. *et al.* A diagnostic host response biosignature for COVID-19 from RNA profiling of
330 nasal swabs and blood. *Sci. Adv.* **7**, eabe5984 (2021).
- 331 14. Olsen, S. J. *et al.* Changes in Influenza and Other Respiratory Virus Activity During the
332 COVID-19 Pandemic — United States, 2020–2021. *MMWR Morb. Mortal. Wkly. Rep.* **70**,
333 1013–1019 (2021).
- 334 15. Vanaerschot, M. *et al.* Identification of a Polymorphism in the N Gene of SARS-CoV-2 That
335 Adversely Impacts Detection by Reverse Transcription-PCR. *J Clin Microbiol* **59**, (2020).
- 336 16. World Health Organization. The top 10 causes of death. (2020).
- 337 17. Chow, E. J. & Mermel, L. A. Hospital-Acquired Respiratory Viral Infections: Incidence,
338 Morbidity, and Mortality in Pediatric and Adult Patients. *Open Forum Infect Dis* **4**, ofx006
339 (2017).
- 340 18. Kalantar, K. L. *et al.* IDseq—An open source cloud-based pipeline and analysis service for
341 metagenomic pathogen detection and monitoring. *Gigascience* **9**, (2020).
- 342 19. Singanayagam, A. *et al.* Duration of infectiousness and correlation with RT-PCR cycle
343 threshold values in cases of COVID-19, England, January to May 2020. *Eurosurveillance* **25**,
344 (2020).
- 345 20. Bullard, J. *et al.* Predicting Infectious Severe Acute Respiratory Syndrome Coronavirus 2
346 From Diagnostic Samples. *Clinical Infectious Diseases* **71**, 2663–2666 (2020).
- 347 21. La Scola, B. *et al.* Viral RNA load as determined by cell culture as a management tool for
348 discharge of SARS-CoV-2 patients from infectious disease wards. *Eur J Clin Microbiol Infect*
349 *Dis* **39**, 1059–1061 (2020).
- 350 22. Bray, N. L., Pimentel, H., Melsted, P. & Pachter, L. Near-optimal probabilistic RNA-seq
351 quantification. *Nature Biotechnology* **34**, 525 (2016).
- 352 23. Sonesson, C., Love, M. I. & Robinson, M. D. Differential analyses for RNA-seq: transcript-level
353 estimates improve gene-level inferences. *F1000Research* **4**, 1521 (2015).
- 354

355 **Data Availability**

356 Gene counts for all UCSF samples have been deposited under NCBI GEO accession
357 GSE188678. The New York dataset can be obtained according to the Data Availability statement
358 in the original publication¹². Code for RNA-seq and qPCR SVM classifier development and
359 validation is available at: <https://github.com/czbiohub/Covid-Host-Classifer-Code>.

360 **Table 1. Performance of 2-gene SVM classifiers on RNA-seq data.** Performance of binary
 361 SVM classifiers for predicting COVID-19 status, measured by the area under the curve (AUC).
 362 Where multiple cross-validation (CV) rounds were performed, values indicate mean and standard
 363 deviation. **a)** Performance of indicated 2-gene combinations in the UCSF training (70%) and
 364 testing (30%) sets, **b)** Performance in the external New York dataset.

365
 366

a

	70% Training Set (n=222, 5-fold CV, 10,000 rounds)	30% Testing Set (n=96, 5-fold CV, 10,000 rounds)	30% Testing Set (n=96, trained on 70% training set)
<i>IFI6, GRINA</i>	0.959 (0.005)	0.936 (0.012)	0.934
<i>IFI6, C15orf48</i>	0.949 (0.005)	0.916 (0.012)	0.908
<i>IFI6, GBP5</i>	0.948 (0.005)	0.917 (0.013)	0.905
<i>IFI44L, GBP5</i>	0.944 (0.004)	0.897 (0.014)	0.883
<i>IFI44L, PTAFR</i>	0.934 (0.006)	0.908 (0.015)	0.910
<i>IFI44L, FCGR1A</i>	0.932 (0.004)	0.864 (0.017)	0.859
<i>HERC6, TNIP3</i>	0.923 (0.005)	0.869 (0.013)	0.844
<i>HERC6, GBP5</i>	0.917 (0.005)	0.826 (0.018)	0.841
<i>HERC6, COA3</i>	0.914 (0.005)	0.787 (0.022)	0.816

367
 368

b

	External Dataset (n=553, 5-fold CV, 10,000 rounds)	External Dataset (n=553, trained on 70% UCSF training set)
<i>IFI6, GRINA</i>	0.875 (0.004)	0.883
<i>IFI6, C15orf48</i>	0.894 (0.004)	0.861
<i>IFI6, GBP5</i>	0.902 (0.004)	0.910
<i>IFI44L, GBP5</i>	0.910 (0.004)	0.919
<i>IFI44L, PTAFR</i>	0.897 (0.004)	0.894
<i>IFI44L, FCGR1A</i>	0.898 (0.003)	0.896
<i>HERC6, TNIP3</i>	0.874 (0.004)	0.852
<i>HERC6, GBP5</i>	0.872 (0.005)	0.866
<i>HERC6, COA3</i>	0.851 (0.004)	0.797

369 **Supplementary Tables**

370 **Supplementary Table 1.** Cohort details.

371

Cohort	Description	COVID-19 (n)	Non-Viral ARI (n)	Other Viral ARI (n)	Reference
1. UCSF (RNA-seq, n=318)	Patients tested for COVID-19, CA, 2020	90	169	59	¹⁰ + this study
2. New York (RNA-seq, n=553)	Patients tested for COVID-19, NY, 2020	166	308	79	¹²
3. UCSF (qPCR, n=144)	Patients tested for COVID-19, CA, 2020	72	72		This study
4. UCSF SARS-CoV-2 N-gene variant (qPCR, n=4)	Patients with SARS-CoV-2 N-gene variant, CA, 2020	4	-		This study
5. UCSF SARS-CoV-2 Delta variant (qPCR, n=7)	Patients with SARS-CoV-2 Delta variant, CA, 2021	7	-		This study
6. UCSF SARS-CoV-2 Omicron variant (qPCR, n=3)	Patients with SARS-CoV-2 Omicron variant, CA, 2021	3	-		This study

372

373 **Supplementary Table 2.** Taqman Gene Expression assay IDs for genes tested by RT-qPCR.

374

Gene	Assay ID
<i>IFI6</i>	Hs00242571_m1
<i>GBP5</i>	Hs00369472_m1
<i>RPP30</i>	Hs01124518_m1
<i>IFI44</i>	Hs00197427_m1
<i>IFI44L</i>	Hs00915292_m1
<i>IFI27</i>	Hs01086373_g1
<i>CCL3</i>	Hs00234142_m1
<i>CD274</i>	Hs00204257_m1
<i>GSTA2</i>	Hs00747232_mH

375 **Supplementary Data Files**

376 **Supplementary Data File 1.** Sample metadata for the UCSF and New York RNA-seq cohorts.

377

378 **Supplementary Data File 2.** AUC values and expression fold-changes in the UCSF cohort and
379 in the New York cohort for every possible 'first' gene and for every possible 'second' gene
380 combined with *IFI6* (related to **Figure 1d,e**).

381

382 **Supplementary Data File 3.** ΔC_t values for host genes in all the samples used in qPCR assays
383 (related to **Figure 2a,c,e,f**). Clinical RT-PCR SARS-CoV-2 C_t values are also provided for the
384 COVID-19 positive samples.

385 Figure Legends

386
387 **Figure 1. Development of 2-gene host-based SVM COVID-19 diagnostic classifiers from**
388 **RNA-seq data. a)** Schematic of the greedy feature selection algorithm used to identify top
389 performing 2-gene combinations. **b)** Receiver operating characteristic (ROC) curve
390 demonstrating performance of SVM classifiers using the indicated 2-gene combinations. The
391 classifiers were trained on the UCSF training set and applied to the UCSF testing set. AUC = area
392 under the ROC curve. **c)** Expression scatter plots and distributions of the representative 'first' and
393 'second' genes *IFI6* and *GBP5*, respectively, in the full UCSF cohort. Shown are variance-
394 stabilized gene expression values following standardization. Color indicates patient group.
395 **d)** Scatter plot of the AUC of COVID-19 diagnostic classifiers that rely on each single gene,
396 calculated using 5-fold cross-validation (y-axis), against \log_2 fold-change (\log_2FC) of the gene
397 between the COVID-19 and non-viral samples (x-axis). Both metrics were averaged between the
398 full UCSF cohort and the New York cohort. **e)** Scatter plot of the AUC of COVID-19 diagnostic
399 classifiers that rely on the combination of *IFI6* and each possible 'second' gene, calculated using
400 5-fold cross-validation (y-axis), against \log_2 fold-change (\log_2FC) of the 'second' gene between
401 the COVID-19 and other viral samples (x-axis). Both metrics were averaged between the full
402 UCSF cohort and the New York cohort. The AUC of an *IFI6*-only classifier is shown for reference.

403
404
405
406 **Figure 2. Performance of 2-gene SVM COVID-19 diagnostic classifiers in qPCR assays.**
407 **a)** Expression differences determined by qPCR in a new cohort of patients with (n=72) or without
408 (n=72) COVID-19 for several 'first' ISGs selected for their predictive value and fold-change in the
409 RNA-seq data. Shown are boxplots of ΔCt values, using *RPP30* as the reference gene,
410 normalized to the median of the COVID-19 group. Statistical significance was assessed using a
411 one-sided Mann-Whitney test with Bonferroni correction. **b)** ROC curves demonstrating
412 performance of SVM classifiers relying on single ISGs for distinguishing the samples with or
413 without COVID-19 using the qPCR data, estimated by 5-fold cross-validation. **c)** Expression
414 differences determined by qPCR between the samples with COVID-19 (n=72) and samples with
415 other viral ARIs (n=17; a subset of the samples from Figure 1) for several 'second' genes selected
416 for their predictive value and fold-change in the RNA-seq data. Shown are boxplots of ΔCt values,
417 using *RPP30* as the reference gene, and normalized to the median of the COVID-19 group.
418 Statistical significance was assessed using a one-sided Mann-Whitney test with Bonferroni
419 correction. **d)** ROC curves demonstrating performance of SVM classifiers relying on 2-gene
420 combinations for distinguishing samples with COVID-19 from samples with other viral ARIs using
421 the qPCR data, estimated by 5-fold cross-validation. **e)** Probability of COVID-19 predicted by the
422 *IFI6+GBP5* classifier using the qPCR data for samples with the Omicron variant (n=3), the N-
423 gene variant (n=4), and the Delta variant (n=7). The classifier was trained on the samples with
424 and without COVID-19 shown in a). **f)** Probability of COVID-19 predicted by the *IFI6+GBP5*
425 classifier using the qPCR data for n=7 samples without COVID-19 before and after trace
426 contamination from a sample with high SARS-CoV-2 viral load. The classifier was trained on the
427 samples with and without COVID-19 shown in a). Statistical significance was assessed using a
428 one-sided, paired Mann-Whitney test.
429 ns = $P > 0.05$, * = $P < 0.05$, ** = $P < 0.01$, *** = $P < 0.001$, **** = $P < 0.0001$.

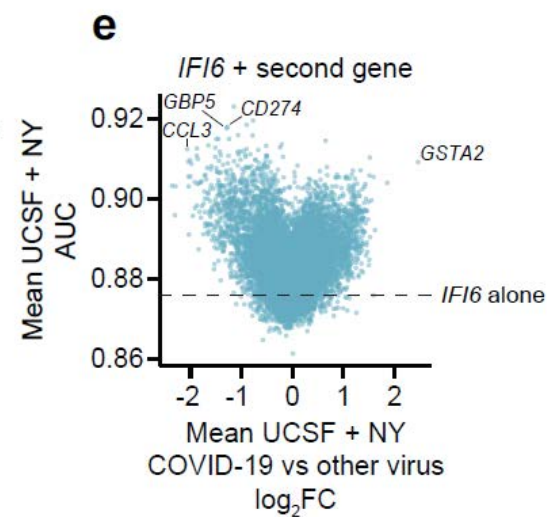
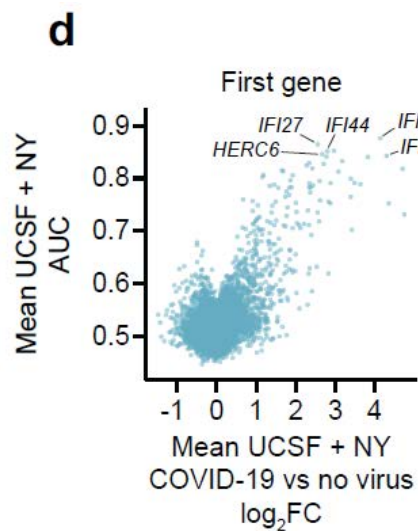
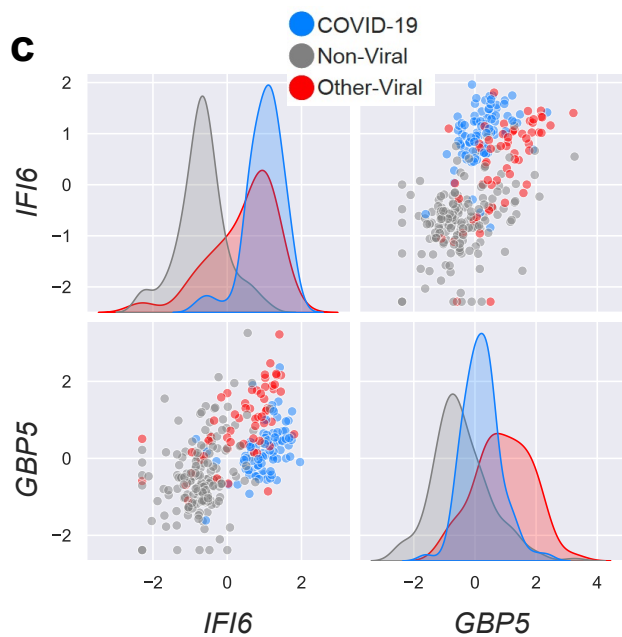
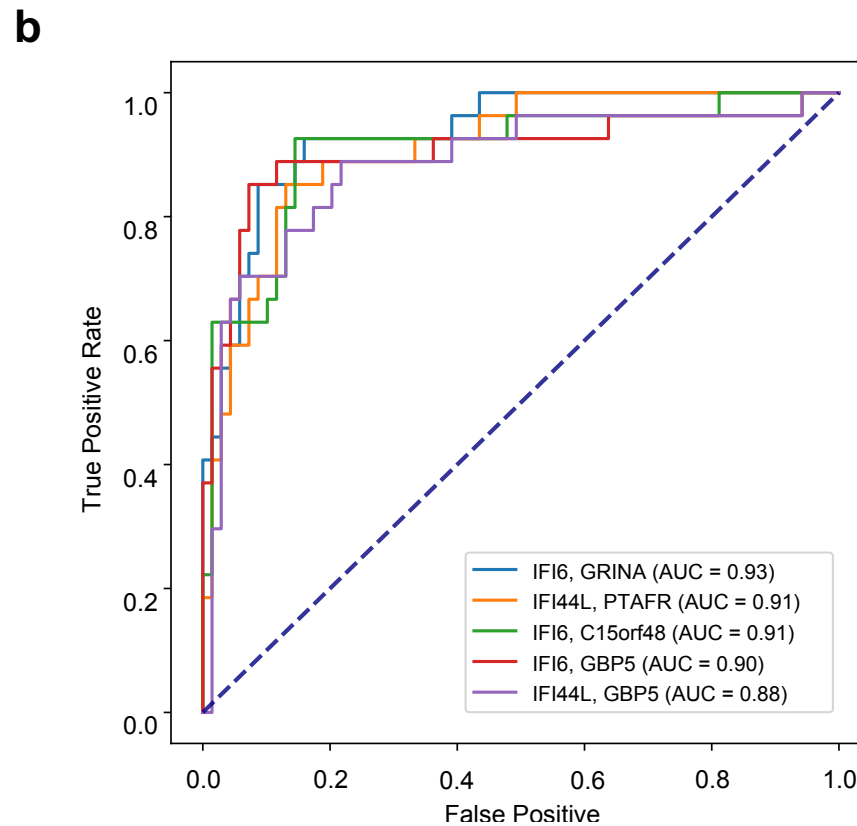
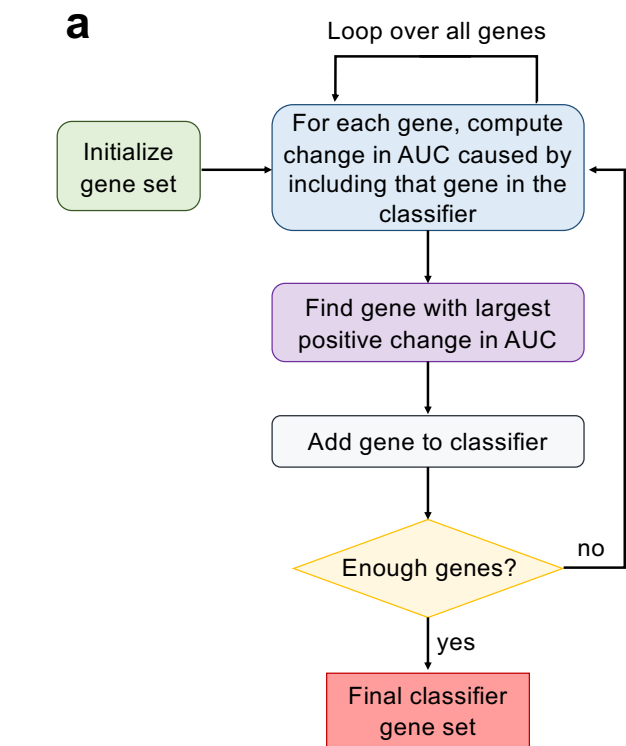


Figure 1

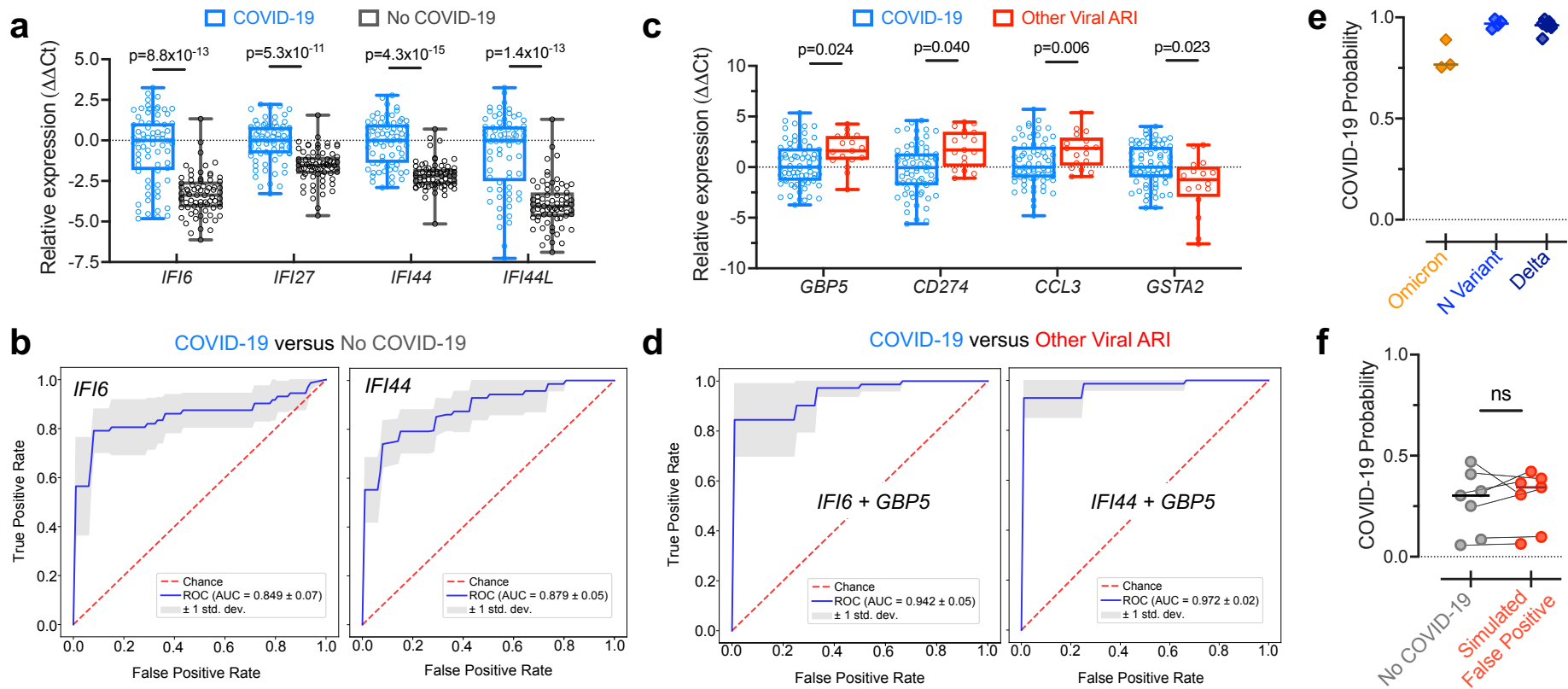


Figure 2